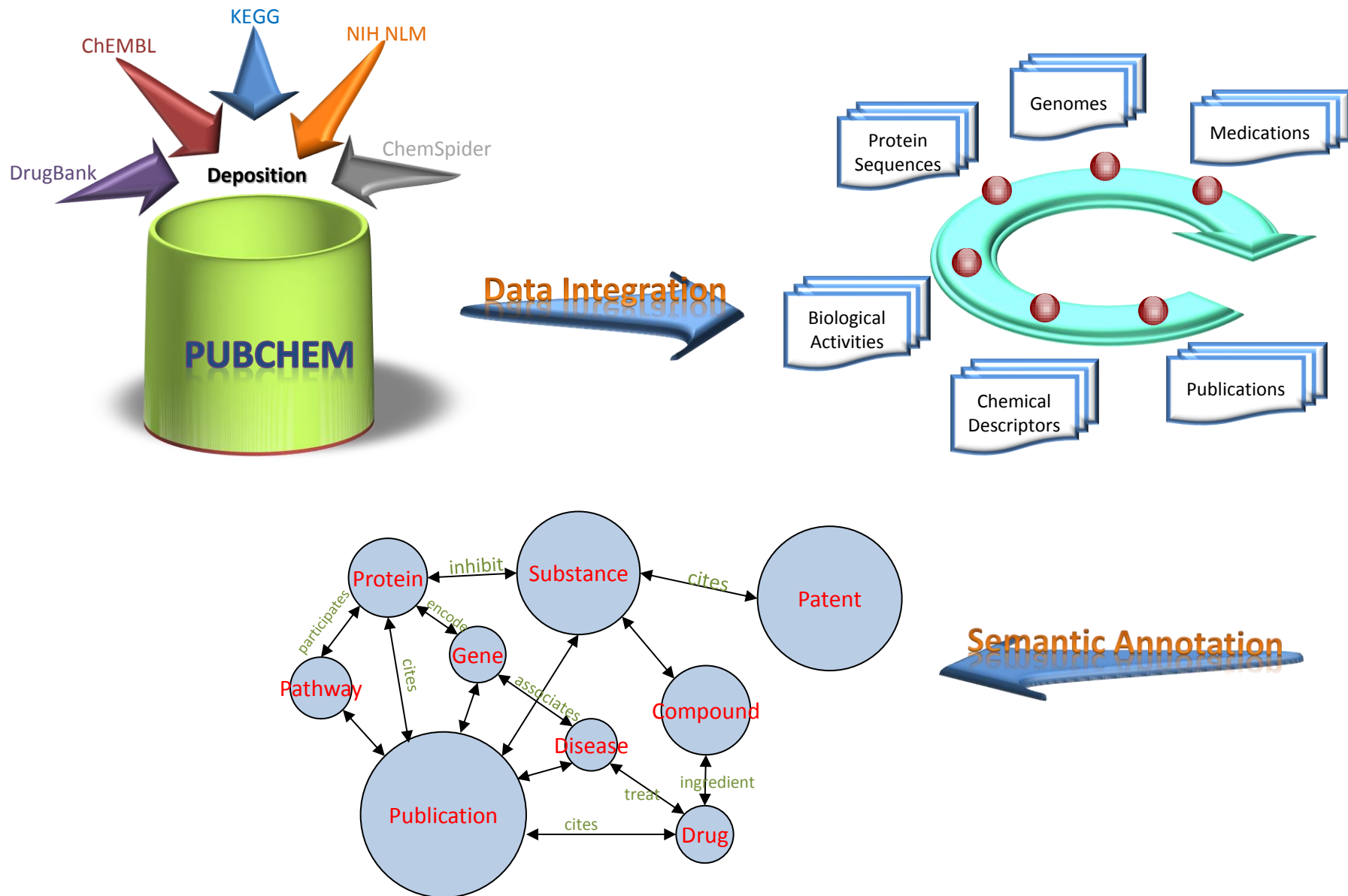# *Semantic Annotation of PubChem Database*

Gang Fu, Colin Batchelor, Michel Dumontier, Janna Hastings, Hande Küçük, Stephan Schurer, Uma Vempati, Egon Willighagen, Evan Bolton

**Pub<span>O</span>hem**

# *OUTLINE*

➢ Background Information

➢ Ontology-based Data Integration

➢ PubChem RDF URI References

➢ PubChem RDF Subdomains

➢ REST Interface

➢ Use Cases

# Ontology-based Data Integration

**Pub Chem**

|  | Standardized Ontologies | |
|---|---|---|
| Prefix | Namespace | Vocabularies |
| rdfs | http://www.w3.org/2000/01/rdf-schema# | RDF Schema |
| rdf | http://www.w3.org/1999/02/22-rdf-syntax-ns# | RDF |
| owl | http://www.w3.org/2002/07/owl# | OWL |
| xsd | http://www.w3.org/2001/XMLSchema# | XML Schema |
| chebi | http://purl.obolibrary.org/obo/ | ChEBI and |
| uo | | UO |
| sio | http://semanticscience.org/resource/ | CHEMINF and |
| cheminf | | SIO |
| skos | http://www.w3.org/2004/02/skos/core# | SKOS |
| obo | http://purl.obolibrary.org/obo/ | BFO and OBI and IAO |
| bao | http://www.bioassayontology.org/bao# | BAO |
| qudt | http://data.nasa.gov/qudt/owl/qudt# | QUDT |
| cito | http://purl.org/spar/cito/ | CiTO |
| fabio | http://purl.org/spar/fabio/ | FaBio |
| ops | http://www.openphacts.org/units/ | Open PHACTS vocabulary |
| pr | http://purl.obolibrary.org/obo/pr# | PRO |
| go | http://purl.obolibrary.org/obo/go# | GO |
| dcterms | http://purl.org/dc/terms/ | DCMI Terms |
| pav | http://purl.org/pav/ | PAV |
| foaf | http://xmlns.com/foaf/0.1/ | FOAF vocabulary |

# PubChem RDF Subdomains

## PubChemRDF Subdomains

| Prefix | Namespace |
|---|---|
| compound | http://rdf.ncbi.nlm.nih.gov/pubchem/compound/ |
| substance | http://rdf.ncbi.nlm.nih.gov/pubchem/substance/ |
| descr | http://rdf.ncbi.nlm.nih.gov/pubchem/descriptor/ |
| syno | http://rdf.ncbi.nlm.nih.gov/pubchem/synonym/ |
| bioassay | http://rdf.ncbi.nlm.nih.gov/pubchem/bioassay/ |
| measuregroup | http://rdf.ncbi.nlm.nih.gov/pubchem/measuregroup/ |
| endpoint | http://rdf.ncbi.nlm.nih.gov/pubchem/endpoint/ |
| protein | http://rdf.ncbi.nlm.nih.gov/pubchem/protein/ |
| domain | http://rdf.ncbi.nlm.nih.gov/pubchem/domain/ |
| biosystem | http://rdf.ncbi.nlm.nih.gov/pubchem/biosystem/ |
| gene | http://rdf.ncbi.nlm.nih.gov/pubchem/gene/ |
| reference | http://rdf.ncbi.nlm.nih.gov/pubchem/reference/ |
| nbr | http://rdf.ncbi.nlm.nih.gov/pubchem/neighbor/ |
| source | http://rdf.ncbi.nlm.nih.gov/pubchem/source/ |
| vocab | http://rdf.ncbi.nlm.nih.gov/pubchem/vocabulary# |

## PubChem RDF URI References

http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID60823

http://rdf.ncbi.nlm.nih.gov/pubchem/substance/SID103554720

http://rdf.ncbi.nlm.nih.gov/pubchem/descriptor/CID60823_Molecular_Weight

http://rdf.ncbi.nlm.nih.gov/pubchem/synonym/MD5_9a05646d461669f86de312d88ab5748a

http://rdf.ncbi.nlm.nih.gov/pubchem/bioassay/AID1788

http://rdf.ncbi.nlm.nih.gov/pubchem/measuregroup/AID447528
http://rdf.ncbi.nlm.nih.gov/pubchem/measuregroup/AID1788_1
http://rdf.ncbi.nlm.nih.gov/pubchem/measuregroup/AID363_PMID16161995

http://rdf.ncbi.nlm.nih.gov/pubchem/endpoint/SID103164874_AID443491
http://rdf.ncbi.nlm.nih.gov/pubchem/endpoint/SID99445338_AID2202_1
http://rdf.ncbi.nlm.nih.gov/pubchem/endpoint/SID8033500_AID363_PMID10395478

http://rdf.ncbi.nlm.nih.gov/pubchem/protein/GI124375976

http://rdf.ncbi.nlm.nih.gov/pubchem/reference/PMID10395478

http://rdf.ncbi.nlm.nih.gov/pubchem/source/ChEMBL
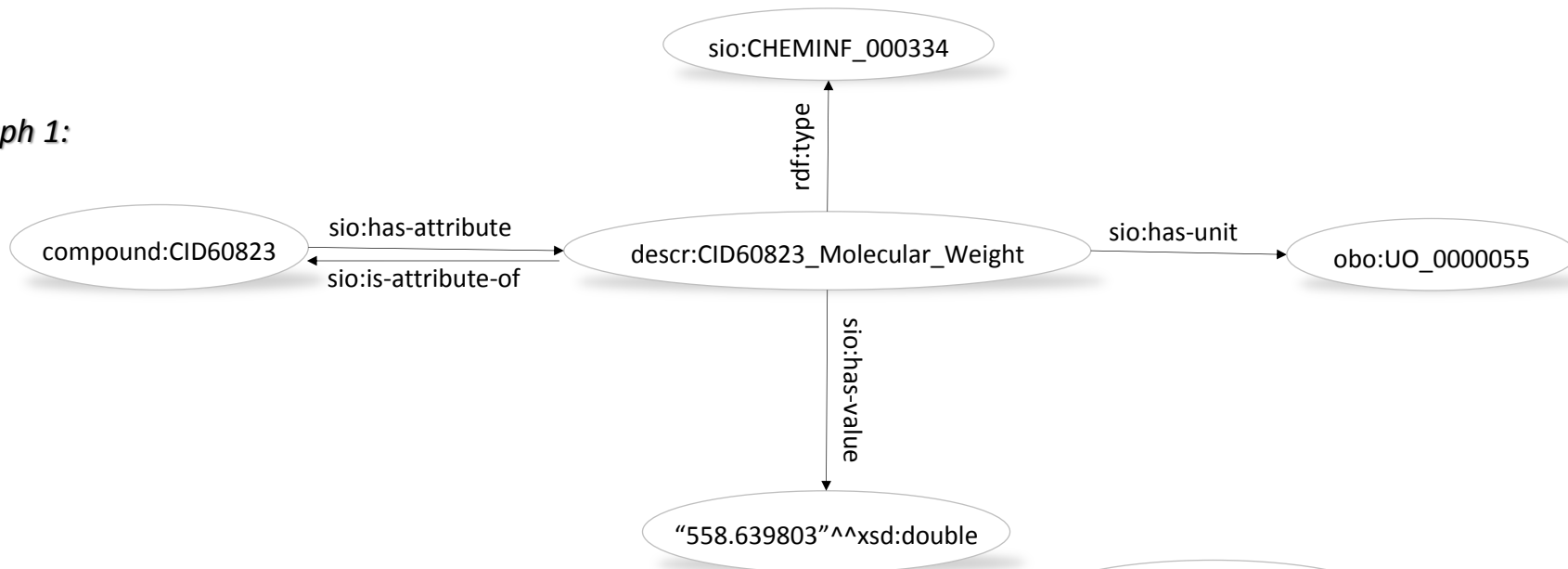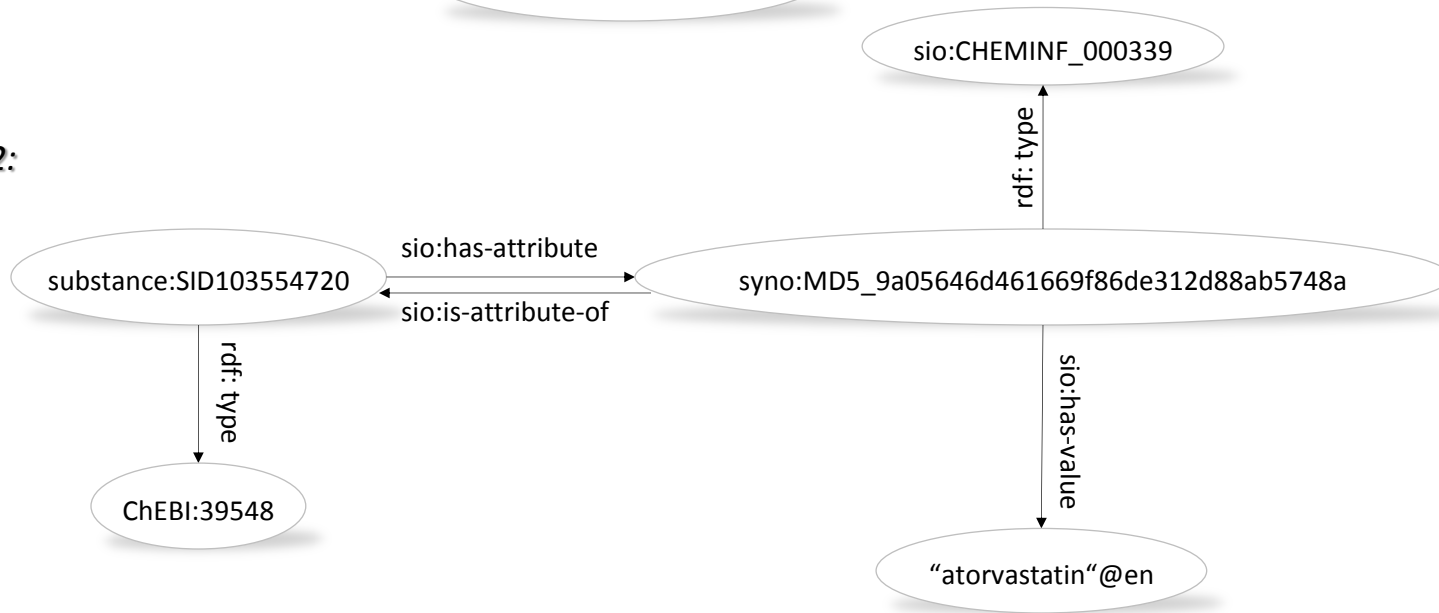
http://rdf.ncbi.nlm.nih.gov/pubchem/neighbor/CID60823_CID68019409_2DSimilarity

http://rdf.ncbi.nlm.nih.gov/pubchem/neighbor/CID60823_CID11330946_3DSimilarity

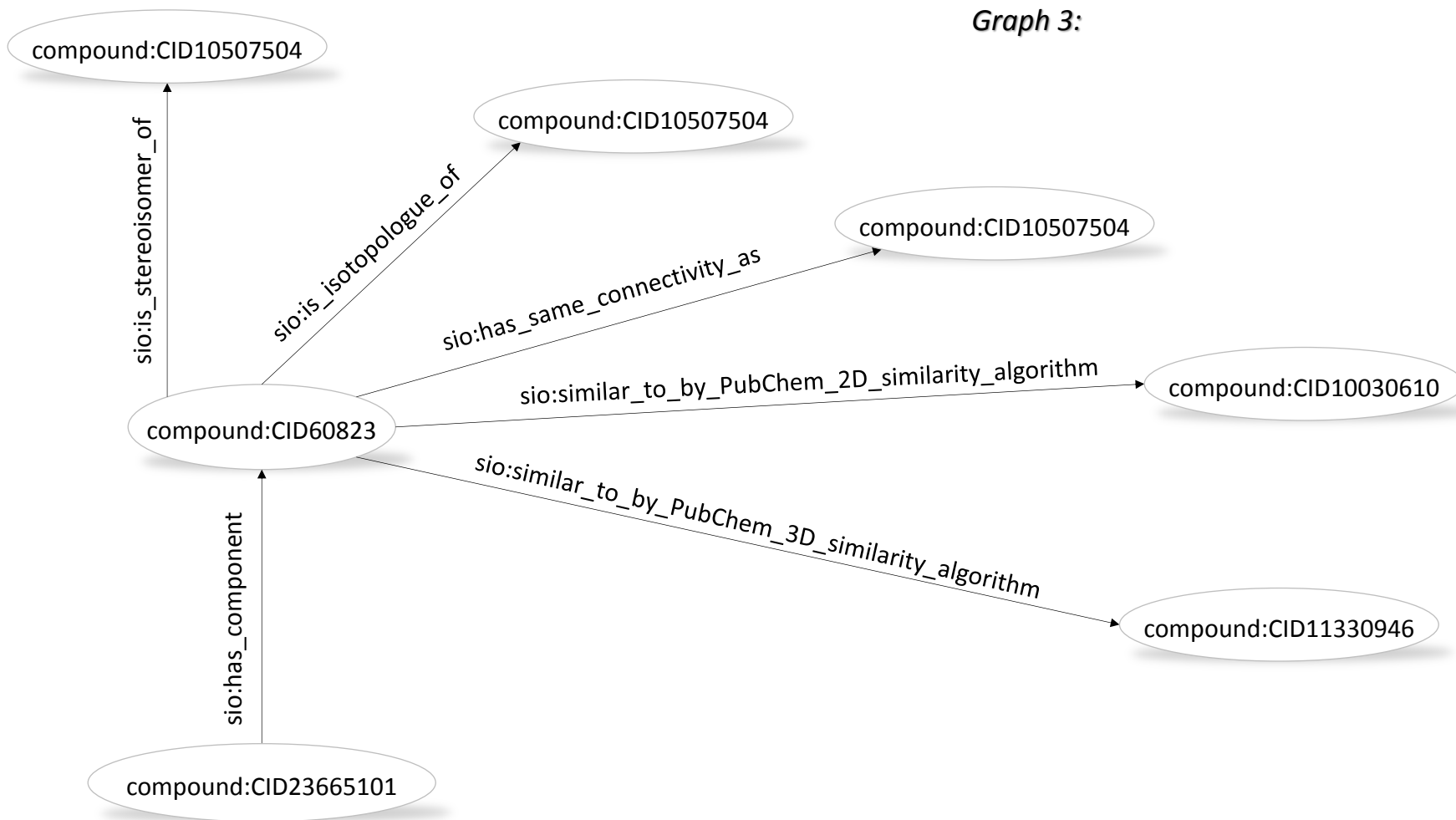http://rdf.ncbi.nlm.nih.gov/pubchem/neighbor/GI548481_GI129900_SequenceSimilarity

*Graph 1:*

sio:CHEMINF_000334

rdf:type

compound:CID60823 —sio:has-attribute→ descr:CID60823_Molecular_Weight ←sio:is-attribute-of—

descr:CID60823_Molecular_Weight —sio:has-unit→ obo:UO_0000055

sio:has-value

"558.639803"^^xsd:double

*Graph 2:*

sio:CHEMINF_000339

rdf: type

substance:SID103554720 —sio:has-attribute→ syno:MD5_9a05646d461669f86de312d88ab5748a ←sio:is-attribute-of—

rdf: type

ChEBI:39548

sio:has-value

"atorvastatin"@en

*Graph 3:*

compound:CID10507504

compound:CID10507504

compound:CID10507504

sio:is_stereoisomer_of

sio:is_isotopologue_of

sio:has_same_connectivity_as

sio:similar_to_by_PubChem_2D_similarity_algorithm

compound:CID60823

compound:CID10030610

sio:similar_to_by_PubChem_3D_similarity_algorithm

compound:CID11330946

sio:has_component

compound:CID23665101

*Graph 4:*



compound:CID60823

compound:CID10030610

*sio:refers-to*

*sio:refers-to*

nbr:CID60823_CID10030610_2DSimilarity

rdf: type

vocab:2D_structural_similarity

sio:has-measurement-value

sio:is-output-of

nbr:CID60823_CID10030610_2DTanimotoScore

rdf: type

vocab:2D_Fingerprint_TanimotoScore

sio:CHEMINF_000333

sio:has-value

"0.98"^^xsd:double

Pub**C**hem

*Graph 5:*

compound:CID60823

compound:CID11330946

nbr:CID60823_CID11330946_3DSimilarity

*sio:refers-to*

*sio:refers-to*

*rdf: type*

vocab:3D_structural_similarity

*sio:has-measurement-value*

*sio:has-measurement-value*

vocab:3D_Shape_
TanimotoScore

nbr:CID60823_CID11330946
_3DShapeTanimotoScore

*rdf: type*

vocab:3D_Feature_
TanimotoScore

*rdf: type*

nbr:CID60823_CID11330946
_3DFeatureTanimotoScore

*sio:has-value*

*sio:is-output-of*

*sio:is-output-of*

*sio:has-value*

sio:CHEMINF_000333

"0.88"^^xsd:double

"0.59"^^xsd:double

*Graph 6:*

*Graph 7:*

*Graph 8:*

# PubChem RDF Linkage between Subdomains

# PubChem RDF REST Interface

- http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID2244.rdf

- http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID2244.xml

- http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID2244.rdfxml

- http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID2244.html

- http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID2244.turtle

- http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID2244.ttl

- http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID2244.json

- http://rdf.ncbi.nlm.nih.gov/pubchem/compound/CID2244.ntriples

| MIME Type | HTTP Accept Header or URI extension |
|---|---|
| application/rdf+xml+abbrev | default |
| | application/rdf+xml+abbrev |
| | rdfxml-abbrev |
| application/rdf+xml | application/rdf+xml |
| | text/rdf |
| | rdfxml |
| | rdf |
| | xml |
| application/xhtml+xml | application/xhtml+xml |
| | text/html |
| | html |
| application/x-turtle[a] | application/rdf+n3 |
| | application/turtle |
| | application/x-turtle |
| | turtle |
| | ttl |
| application/json[b] | application/json |
| | text/json |
| | json |
| text/plain | text/plain |
| | ntriples |
| text/rdf+n3 | text/n3 |
| | text/rdf+n3 |
| | Accept: n3 |

| HTTP Status | Description |
|---|---|
| 400 | Bad Query URL or Request URI |
| 404 | Input URI is invalid or cannot be identified in databases |
| 405 | MIME output format is unspecified or invalid |
| 500 | Some problem on the server side occurs |
| 504 | The request timed out (over 28 second) |

# *PubChem RDF Utility*

**Pub⬡Chem**



API functions

Applications

SPARQL queries

HTTP protocols

SPARQL endpoint

Jena core library

Jena ARQ library
Virtuoso isql

Fuseki HTTP service
Virtuoso conductor HTTP service

**Ontology-based Data Integration**
*Rule-based reasoning, forward and backward chaining, consistency validation ...*

In memory
Storage

Jena TDB Library

Virtuoso Jena Provider
JDBC

JDBC
R2RML (RDB2RDF)

Jena TDB Triple Store

Virtuoso RDF Quad Store

RDBMS

Single
Server
can
handle
millions
of triples

Cluster of
Servers
can
handle
billions of
triples

Persistent
Database
Storage

PubChemRDF CGI program:
URI-resolving REST interface

PubChem Database RDFization:
FTP dump files in Turtle representation

# PubChem RDF Utility

**Pub**C**hem**

- ➢ PubChem RDF is intended for ontology-based data integration

- ➢ PubChem databases have been semantically exposed to linked open data

- ➢ REST interface can be accessed to resolve URI references

- ➢ FTP dump files can be bulk-loaded into open source triples stores

*Acknowledgement*

**Pub⬡hem**

**NCBI Structure Group:**
Steve Bryant
Evan Bolton
Yanli Wang
Paul Thiessen
Siqian He
Jeff Zhang
Tugba Suzek
Lianyi Han
Jane He
Jiyao Wang
Tiejun Cheng

**External Collaborators:**
Colin Batchelor
Michel Dumontier
Janna Hastings
Hande Küçük
Stephan Schurer
Uma Vempati
Egon Willighagen

# Thank you and Questions!